

雑談対話モデルの関連性向上のための 関連語を優先した応答文生成手法の検討

Study on response generation prioritizing related words towards relevance improvement for chat dialogue model

麻生大聖^{1*} 滝口哲也¹ 有木康雄¹
Aso Taisei¹ Takiguchi Tetsuya¹ Ariki Yasuo¹

¹ 神戸大学システム情報学研究科

¹ Graduate School of System Informatics, Kobe University, Department of Information Science

Abstract: The purpose of this study is to improve the relevance of responses generated by the already learned chat dialogue model. Neural chat dialogue models may generate dull responses (e.g. "I'm not sure.") or non-related responses to user utterances. To reduce such responses, in this study, we try to enhance the output probabilities of the words related to user utterances in the knowledge graph during inference. We used several decoding strategies (Greedy, Beam Search, Sampling, Top-k Sampling, Top-p Sampling, MMI-antiLM) and analyzed the performance improvement when the proposed method is applied to each strategy. We evaluated responses in terms of diversity, appropriateness, and relevance.

1 はじめに

近年、IoT化に伴って会話型インターフェースが拡大しており、雑談を行う機能が人間とのやりとりを円滑にする重要な役割を担い始めている。しかし、ニューラル雑談モデルでは無難でつまらない応答や、ユーザの発話と関連性が低い応答をする場合がある。

そこで我々は、推論時にユーザの発話語に関連する単語を知識グラフから抽出し、自己相互情報量 (Pointwise Mutual Information; PMI) に基づいたバイアスをその関連語のモデル出力確率にかけることによってユーザの発話に対する雑談応答の関連性を向上させる手法を提案する。

提案手法は、モデルを再学習することなく雑談応答の品質を調節でき、またモデル依存性が低い点も利点である。生成文はデコード戦略によって大きく変化するため、Beam Searchなどの複数のデコード戦略に対して提案手法を適用し比較する。

2 ConceptNet

本研究では、ユーザの発話語に関連する単語を抽出するために、知識グラフである ConceptNet[1] を使用

する。ConceptNetとは、自然言語の単語およびフレーズ（それらをエンティティと呼ぶ）が様々な述語によって関連付けられた大規模な意味グラフであり、「パリはフランスの首都」のような事実だけでなく、「犬はペット」のような日常的な関係も含まれる点が特徴的である[2]。ConceptNetにおいて任意の二単語は、

music $\xrightarrow{\text{RelatedTo}}$ song $\xrightarrow{\text{CreatedBy}}$ songwriter

のように1つ以上の述語によって結ばれる可能性がある。本研究では、二単語間を結ぶために必要な述語の数をそれらの距離と定義し、距離が一定以下である二単語は関連性が高いと見なす。

3 従来のデコード戦略

本研究では、Seq2seqモデルを雑談応答生成に用いることを想定している。ある時刻 t の Decoder の最終層の出力ベクトル $O^t = \{o_1^t, \dots, o_K^t\}$ が与えられたときに、以下の6つのデコード戦略に基づいて単語を出力することを毎時刻繰り返す、応答文を生成する。 K はモデルの語彙数を表し、 O^t は語彙集合 $V = \{v_1, \dots, v_K\}$ に対応する。出力単語を y_t とし、 $y_t = v_{\hat{y}_t}$ が成り立つように \hat{y}_t を定義する。

Greedy: (1) 式でモデル出力が最大となる単語を逐次的に選択する。

$$\hat{y}_t = \operatorname{argmax} O^t \quad (1)$$

*連絡先：神戸大学システム情報学研究科
〒657-8501 兵庫県神戸市灘区六甲台町 1-1
自然科学総合研究棟 3 号館 805 号室
E-mail: taisei.aso@stu.kobe-u.ac.jp

Beam Search (BS): (2) 式で系列尤度が上位 B 以内の単語系列集合 $Y_{\leq t}$ を保持する。 \mathcal{Y}_t は時刻 t における探索単語系列 $y_{\leq t}^b$ の集合である。 取り出した B 個の系列の尤度が上位 B 以内であることを表現するために総和を取っている。 単語の生起確率は、(3) 式で softmax 関数により計算する。

$$Y_{\leq t} = \operatorname{argmax}_{y_{\leq t}^1, \dots, y_{\leq t}^B \in \mathcal{Y}_t} \sum_{b=1}^B \left(\prod_{t'=1}^t Pr(y_{t'}^b) \right) \quad (2)$$

$$Pr(y_t) = \operatorname{softmax}(O_{\hat{y}_t}^t) = e^{o_{\hat{y}_t}^t} / \sum_{k=1}^K e^{o_k^t} \quad (3)$$

Sampling[3]: (4) 式で求めた単語の生起確率分布を重みとして、単語 $\hat{y}_t \sim tPr(y_t)$ をサンプリングする。 tmp は softmax 関数の温度パラメータであり、モデル出力分布の滑らかさを調節できる。

$$tPr(y_t) = \operatorname{softmax} \left(\frac{O_{\hat{y}_t}^t}{tmp} \right) = e^{\frac{o_{\hat{y}_t}^t}{tmp}} / \sum_{k=1}^K e^{\frac{o_k^t}{tmp}} \quad (4)$$

Top-k Sampling[4]: モデル出力 O^t から上位 $topk$ 個の候補単語を取り出し、それらに絞って (4) 式の温度付き softmax 関数で生起確率分布に正規化し、単語を重み付きサンプリングする。

Top-p Sampling[5]: (4) 式で求めた生起確率が上位の単語から累積確率が $topp$ を超えるまで取り出し、それらに絞って合計確率が 1 になるように正規化したのち、単語を重み付きサンプリングする。

MMI-antiLM[6]: (5) 式により、モデル出力 O_t から言語モデルのバイアス出力 U_t をペナルティとして減算することで応答文の多様性を向上させる。 (5) 式は入力文と生成文の相互情報量を最大化する効果がある。 U^t は入力文が与えられない時のモデル出力ベクトルであり、実験的には (Attention ベクトルのための) Encoder の出力ベクトルと Decoder の初期状態を 0 にして求めた。 λ はペナルティの強弱を調整するハイパーパラメータである。 応答文の後半になるほど言語モデルのバイアス出力が強くなり、MMI-antiLM では文構造の破綻が起きる可能性が高いため、最初の 5 単語は (5) 式により生成し、以降は (1) 式により Greedy で生成する。

$$\hat{y}_t = \operatorname{argmax}(O^t - \lambda \cdot U^t) \quad (5)$$

4 提案手法

提案手法では、雑談応答の関連性を向上させるために、ユーザの発話語に関連する単語のモデル出力にバイアスをかける。 バイアスをかける手法と、各デコード戦略への適用方法について述べる。

4.1 モデル出力への関連語バイアス

(6)(7)(8) 式でモデル出力ベクトル O の全ての要素 o_k を $EE(o_k)$ に置換する。 EE (Entity Enhancer) の処理は時刻には依存せず、全ての時刻におけるモデル出力に同じバイアスをかける。 EE の概略図を図 1 に示す。

$$EE(o_k) = \left(1 + \alpha \cdot \max_{x \in X} r(x, v_k) \right) \cdot o_k \quad (6)$$

$$r(x, y) = \begin{cases} \max(0, PMI(x, y)) & d(x, y) \leq 2 \\ 0 & otherwise \end{cases} \quad (7)$$

$$PMI(x, y) = \log_2 \frac{P(x, y)}{P_X(x) \cdot P_Y(y)} \quad (8)$$

(6) 式において X は入力単語系列で、 $r(x, y)$ は二単語間の関連度を表す。 入力単語系列に対する最大の関連度に、ハイパーパラメータ α をかけた値をバイアスとする。 α によって関連語の生成を優先する度合いを調整することができる。

二単語間の関連度は、(7) 式で自己相互情報量 PMI に基づいて 0 以上の値として計算される。 ConceptNet 上での二単語間の距離が $d(x, y)$ と表され、距離が 2 より大きい二単語間の関連度は 0 とする。

二単語間の自己相互情報量は (8) 式で表せる。 $P_X(x)$ は x が入力文に現れる確率、 $P_Y(y)$ は y が応答文に現れる確率、 $P(x, y)$ は x が入力文に現れてかつ y が応答文に現れる確率を表す。 これらの確率は本研究では学習用対話コーパスで計算される。 $P(x, y) = 0$ の場合は $-\infty$ への発散を防ぐため $PMI(x, y) = 0$ とする。 自己相互情報量は、雑談対話中により多く共起する特徴的な (低頻度な) 二単語間ほど大きい値を取る。

$EE(O)$ などベクトルが入力される場合は、全ての要素を EE で置換したベクトルを返すと定義する。

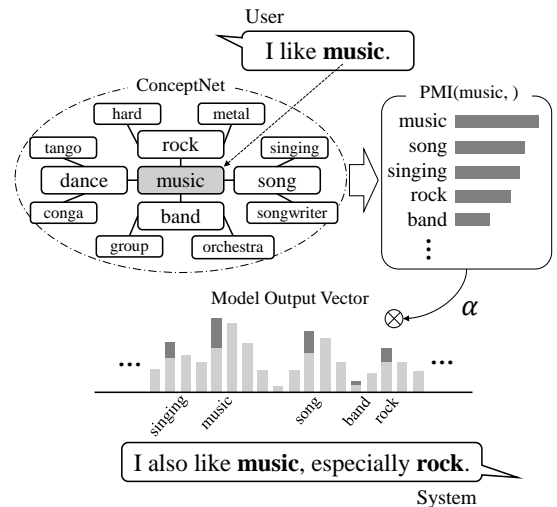


図 1: 提案手法 EE の概略図

4.2 デコード戦略への適用

各デコード戦略への提案手法 EE の適用方法を以下に述べる。

Greedy: (1) 式の代わりに (9) 式で単語を逐次的に選択する。 O^t の代わりに $EE(O^t)$ を用いる。

$$\hat{y}_t = \operatorname{argmax} EE(O^t) \quad (9)$$

Beam Search (BS): (2) 式の代わりに (10)(11) 式を用いて単語系列を探索する。系列尤度を求める際に、現在時刻のモデル出力にのみ EE を適用する。関連語をとにかく多く含む系列を探索するような戦略を避け、生成文の自然性（適切性）を維持するねらいがある。

$$Y_{\leq t} = \operatorname{argmax}_{y_{\leq t}^1, \dots, y_{\leq t}^B \in \mathcal{Y}_t} \sum_{b=1}^B \left(\prod_{t'=1}^{t-1} Pr(y_{t'}^b) \cdot ePr(y_t^b) \right) \quad (10)$$

$$ePr(y_t) = \operatorname{softmax}(EE(O_{\hat{y}_t}^t)) = \frac{e^{EE(o_{\hat{y}_t}^t)}}{\sum_{k=1}^K e^{EE(o_k^t)}} \quad (11)$$

Sampling: O^t の代わりに $EE(O^t)$ を用いる。

Top-k Sampling: O^t の代わりに $EE(O^t)$ を用いる。

Top-p Sampling: O^t の代わりに $EE(O^t)$ を用いる。

MMI-antiLM: 最初の 5 単語は (12) 式により生成し、以降は (9) 式で生成する。

$$\hat{y}_t = \operatorname{argmax}(EE(O^t - \lambda \cdot U^t)) \quad (12)$$

5 実験

5.1 データセット

本研究では、雑談対話コーパスおよび ConceptNet が整形されたデータセットである Commonsense Conversation Dataset[2] を用いる。雑談対話コーパスは Reddit から単ターン対話が収集されており、入力文と応答文のいずれの単語対も ConceptNet 上で距離 1 以内で接続できない場合は除外されている。すなわち、常識的なつながりのある対話のみがフィルタリングされている。Commonsense Conversation Dataset の ConceptNet は簡単のため複数単語を含むエンティティが削除されており、21,471 個のエンティティと 44 種類の述語を含む。(エンティティ, 述語, エンティティ) の三つ組 (トリプル) は 120,850 個存在する。

雑談対話コーパスのうち、50 万対話を学習に、1000 対話を評価に用いる。また、データセットに含まれる ConceptNet は全て用いる。

5.2 実験設定

注意機構付きの GRU Encoder-Decoder モデルを用いた。Encoder は双方向で各方向はそれぞれ 2 層、Decoder は 2 層とした。モデルサイズは各層 128 次元で、埋め込みベクトルの初期値として 300 次元の GloVe を用いた。最適化手法として初期学習率を $5e-4$ とした Adam を使用した。過学習防止のために Dropout を 0.1 に設定し、入力文に対するロバスト性向上のために学習時に入力単語を 5% の確率でランダムな単語に変化させた。ミニバッチサイズを 128 とし、20epoch 学習させた。

各デコード戦略のパラメータは、グリッドサーチにより応答文の多様性と適切性のバランスをみて表 1 に決定した。Beam Search では候補単語系列の尤度を系列長で正規化することで、短文しか生成されなくなることを防いだ。また、全てのデコード戦略において、繰り返し同じ単語を生成することを抑制するために repetition suppressor[7] を全てのモデル出力に適用した。

表 1: デコード戦略のパラメータ

Strategy	Parameters
Beam Search	$B = 5$
Sampling	$tmp = 0.6$
Top-k Sampling	$topk = 32, tmp = 0.6$
Top-p Sampling	$topp = 0.6, tmp = 0.8$
MMI-antiLM	$\lambda = 0.6$

5.3 評価指標

全ての生成文を以下の多様性、適切性、関連性の指標で評価した。

多様性 (Diversity): 全ての生成文に含まれる n-gram のうち異なるものの割合 DIST-n を計算する。DIST-1 と DIST-2 を用いる。

適切性 (Appropriateness): 機械翻訳や要約生成に用いられる BLEU-1, BLEU-2, ROUGE-L, NIST, METEOR を用いる。

関連性 (Relevance): (13) 式で定義される文単位での一貫性指標 PMI を用いる。 X は入力単語系列で、 Y は応答単語系列である。

$$PMI = \frac{1}{|Y|} \cdot \sum_{y \in Y} \max_{x \in X} PMI(x, y) \quad (13)$$

また、入力単語のいずれかとの ConceptNet 上での距離が 2 以下である単語が応答文に含まれる平均個数 ENT を定義し用いる。

表 2: 各デコード戦略および提案手法適用時の客観評価の比較

Strategy	Length	Diversity		Appropriateness					Relevance	
		DIST-1	DIST-2	BLEU-1	BLEU-2	ROUGE-L	NIST	METEOR	PMI	ENT
Greedy	11.533	4.804	15.105	10.001	1.673	11.994	0.332	9.322	0.582	1.396
+ $EE(\alpha = 0.1)$	9.957	12.685	36.943	9.175	1.380	11.770	0.169	9.073	1.619	2.865
BS	12.514	3.340	10.752	11.644	2.083	12.950	0.500	10.578	0.525	1.754
+ $EE(\alpha = 0.1)$	12.017	5.159	15.748	11.349	2.172	13.096	0.438	10.637	0.755	1.970
BS + reranking(PMI)	17.583	4.595	17.590	13.602	2.748	14.238	1.240	12.729	0.738	2.686
+ $EE(\alpha = 0.1)$	19.036	8.200	28.293	14.691	2.991	14.810	1.412	13.619	1.154	4.451
Sampling	15.968	9.331	41.656	12.433	1.761	12.474	0.942	10.942	0.673	2.717
+ $EE(\alpha = 0.1)$	13.536	14.074	55.584	10.871	1.397	11.703	0.650	9.798	1.533	4.046
Top-k Sampling	15.649	8.154	38.405	12.194	1.648	12.447	0.926	10.757	0.665	2.626
+ $EE(\alpha = 0.1)$	14.007	13.208	52.956	11.100	1.395	11.942	0.733	9.901	1.547	4.229
Top-p Sampling	15.534	7.680	34.203	12.266	1.760	12.610	0.906	10.873	0.663	2.547
+ $EE(\alpha = 0.1)$	13.290	13.213	49.536	11.288	1.452	12.097	0.637	10.401	1.596	4.011
MMI-antiLM	11.296	6.144	20.299	11.170	1.499	13.445	0.328	10.278	0.788	2.209
+ $EE(\alpha = 0.1)$	10.161	14.349	46.338	10.224	1.042	12.401	0.198	9.392	1.940	3.880

表 3: 応答文の比較 (A...BS+reranking(PMI), B...BS+reranking(PMI)+ $EE(\alpha = 0.1)$)

Input	Response
yeah , im home back in cold old england aha	A) i 'm not sure if you 're in the same boat .
	B) i 'm not sure if you 're in northern ireland
i think lenovo is one of the best laptop companies , personally .	A) that 's a good point . i 'm not sure if it was n't the case , but there is no need to be an issue for me and my experience
	B) i 'm not sure what you 're talking about . it 's just a hardware device , but the laptop is n't that bad

5.4 結果と考察

各デコード戦略とそれらに提案手法 EE を適用した時の客観評価の比較を表 2 に示す。Beam Search は候補単語系列のうち最大尤度ものを選択する手法 (BS) と、一貫性指標 PMI が最大ものを選択するリランキング手法 (BS+reranking(PMI)) を用いた。

EE を適用することで、全てのデコード戦略において多様性 (Diversity) や関連性 (Relevance: 一貫性と関連語数) が大きく向上した。しかし、Beam Search 以外のデコード戦略では BLEU などの適切性指標が大きく劣化した。 EE により、本来は生起確率が低かった単語が生成され、文構造の維持が困難になったためであると考えられる。対して、Beam Search は適切性が劣化しにくかった ($\alpha = 0.1$ の時には多くの関連性指標が向上した)。現在時刻までの系列尤度も考慮して単語系列を探索するため、関連語を取り込んでも文構造が破綻しにくいと考えられる。

EE の適用による応答文の変化の例を表 3 に示す。適用後の方がより関連した応答文になっていることが確認できた。

以上のことから、複数系列を保持し探索できる Beam Search は提案手法 EE と相性が良いと考えられるが、多様性が Sampling などよりも劣る。今後は Beam Search の候補系列の多様性を向上させる Group Diverse Beam Search[8] などに適用したり、主観評価を行いたい。

謝辞

本研究の一部は、JSPS 科研費 JP17H01995 の助成を受けたものである。

参考文献

- [1] Catherine Havasi et al.: ConceptNet 3: a Flexible, Multilingual Semantic Network for Common Sense Knowledge, In *Recent Advances in Natural Language Processing*, 2007
- [2] Hao Zhou et al.: Commonsense Knowledge Aware Conversation Generation with Graph Attention, In *Proceedings of IJCAI-ECAL*, pp. 4623–4629, 2018
- [3] Iulian V. Serban et al.: Building End-To-End Dialogue Systems Using Generative Hierarchical Neural Network Models, In *Thirtieth AAAI Conference on Artificial Intelligence*, 2016
- [4] Angela Fan et al.: Hierarchical Neural Story Generation, In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, pp. 889–898, 2018
- [5] Ari Holtzman et al.: The Curious Case of Neural Text Degeneration, *arXiv preprint arXiv:1904.09751*, 2019
- [6] Jiwei Li et al.: A Diversity-Promoting Objective Function for Neural Conversation Models, In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 110–119, 2016
- [7] Ryo Nakamura et al.: Another Diversity-Promoting Objective Function for Neural Dialogue Generation, *arXiv preprint arXiv:1811.08100*, 2018
- [8] Ashwin K Vijayakumar et al.: Diverse Beam Search: Decoding Diverse Solutions from Neural Sequence Models, *arXiv preprint arXiv:1610.02424*, 2016